

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Express Mail No.: EL067144491US

In re application of: LAURILA et al.

Serial No.: 0 /

Filed: Herewith

For: METHOD IN SPEECH RECOGNITION AND A SPEECH RECOGNITION DEVICE

Examiner:

Group No.:

Commissioner of Patents and Trademarks
Washington, D.C. 20231

TRANSMITTAL OF CERTIFIED COPY

Attached please find the certified copy of the foreign application from which priority is claimed for this case:

Country : Finland
Application Number : 990078
Filing Date : 18 January 1999

WARNING: "When a document that is required by statute to be certified must be filed, a copy, including a photocopy or facsimile transmission of the certification is not acceptable." 37 CFR 1.4(f) (emphasis added.)

SIGNATURE OF ATTORNEY

Reg. No.: 24,622

Clarence A. Green

Type or print name of attorney

Tel. No.: (203) 259-1800

Perman & Green, LLPP.O. Address425 Post Road, Fairfield, CT 06430

NOTE: The claim to priority need be in no special form and may be made by the attorney or agent if the foreign application is referred to in the oath or declaration as required by § 1.63.

(Transmittal of Certified Copy [5-4])

#2
jc584 U.S. PTO
09/482277
01/13/00

CERTIFICATE

I, Tuulikki Tulivirta, hereby certify that, to the best of my knowledge and belief, the following is a true translation, for which I accept responsibility, of a certified copy of Finnish Patent Application 990078 filed on 18 January 1999.

Tampere, 16 December 1999



Tuulikki Tulivirta

Tuulikki Tulivirta
Certified Translator (Act 1148/88)

Tampereen Patenttitoimisto Oy
Hermiankatu 6
FIN-33720 TAMPERE
Finland

Method in speech recognition and a speech recognition device

- 5 The present method relates to a method in speech recognition as set forth in the preamble of the appended claim 1, a speech recognition device as set forth in the preamble of the appended claim 8, and a speech-controlled wireless communication device as set forth in the preamble of the appended claim 11.
- 10 For facilitating the use of wireless communication devices, speech recognition devices have been developed, whereby a user can utter speech commands which the speech recognition device attempts to recognize and convert to a function corresponding to the speech command, *e.g.* a command to select a telephone number. A problem in the
- 15 implementation of speech control has been for example the fact that different users say the speech commands in different ways: the speech rate can be different between different users, so does the speech volume, voice tone, *etc.* Furthermore, speech recognition is disturbed by a possible background noise, whose interference outdoors and in a
- 20 car can be significant. Background noise makes it difficult to recognize words and to distinguish between different words *e.g.* upon uttering a telephone number.
- 25 Some speech recognition devices apply a recognition method based on a fixed time window. Thus, the user has a predetermined time within which s/he must utter the desired command word. After the expiry of the time window, the speech recognition device attempts to find out which word/command was uttered by the user. However, such a method based on a fixed time window has *e.g.* the disadvantage that all
- 30 the words to be uttered are not equally long; for example, in names, the given name is often clearly shorter than the family name. Thus, after a shorter word, more time will be consumed for the recognition than in the recognition of a longer word. This is inconvenient for the user. Furthermore, the time window must be set according to slower speakers so
- 35 that recognition will not be started until the whole word is uttered. When words are uttered faster, a delay between the uttering and the recognition increases the inconvenient feeling.

Another known speech recognition method is based on patterns formed of speech signals and their comparison. Patterns formed of command words are stored beforehand, or the user may have taught desired words which have been formed into patterns and stored. The speech
5 recognition device compares the stored patterns with feature vectors formed of sounds uttered by the user during the utterance and calculates the probability for the different words (command words) in the vocabulary of the speech recognition device. When the probability for a command word exceeds a predetermined value, the speech recognition
10 device selects this command word as the recognition result. Thus, incorrect recognition results may occur particularly in the case of words in which the beginning resembles phonetically another word in the vocabulary. For example, the user has taught the speech recognition device the words "Mari" and "Marika". When the user is saying the word
15 "Marika", the speech recognition device may make "Mari" as the recognition decision, even though the user may not yet have had time to articulate the end of the word. Such speech recognition devices typically use the so-called Hidden Markov Model (HMM) speech recognition method.

20 U.S. patent 4,870,686 presents a speech recognition method and a speech recognition device, in which the determination of the end of words by the user is based on silence; in other words, the speech recognition device examines if there is a perceivable audio signal or not.
25 A problem in this solution is the fact that a too loud background noise may prevent the detection of pauses, wherein the speech recognition is not successful.

30 It is an aim of the present invention to provide an improved method for detecting pauses in speech and a speech recognition device. The invention is based on the idea that a tone band to be examined is divided into sub-bands, and the power of the signal is examined in each sub-band. If the power of the signal is below a certain limit in a sufficient number of sub-bands for a sufficiently long time, it is deduced that there
35 is a pause in the speech. The method of the present invention is characterized in what will be presented in the characterizing part of the appended claim 1. The speech recognition device according to the present invention is characterized in what will be presented in the char-

acterizing part of the appended claim 8. The wireless communication device of the present invention is characterized in what will be presented in the characterizing part of the appended claim 11.

- 5 The present invention gives significant advantages to the solutions of prior art. By the method of the invention, a more reliable detection of a gap between words can be obtained than by methods of prior art. Thus, the reliability of the speech recognition is improved and the number of incorrect and failed recognitions is reduced. Furthermore, the speech
- 10 recognition device is more flexible with respect to manners of speaking by different users, because the speech commands can be uttered more slowly or faster without an inconvenient delay in the recognition or recognition taking place before an utterance has been completed.
- 15 By the division into sub-bands according to the invention, it is possible to reduce the effect of external interference. Spurious signals *e.g.* in a car have typically a relatively low frequency. In solutions of prior art, the energy contained in the whole frequency range of the signal is utilized in the recognition, wherein signals which are strong but have a narrow
- 20 band width reduce the signal-to-noise ratio to a significant degree. Instead, if the frequency range to be examined is divided into sub-bands according to the invention, the signal-to-noise ratio can be improved significantly in such sub-bands in which the proportion of spurious signals is relatively small, which improves the reliability of the rec-
- 25 ognition.

In the following, the present invention will be described in more detail with reference to the appended drawings, in which

- 30 Fig. 1 is a flow chart illustrating the method according to an advantageous embodiment of the invention,
- Fig. 2 is a reduced flow chart showing the speech recognition device according to an advantageous embodiment of the
- 35 invention,

Fig. 3 is a state machine chart illustrating rank-order filtering to be applied in the method according to an advantageous embodiment of the invention, and

5 Fig. 4 is a flow chart illustrating the logic for deducing a pause to be applied in the method according to an advantageous embodiment of the invention.

10 The following is a description on the function of the method according to an advantageous embodiment of the invention, with reference to the flow chart of Fig. 1 and using as an example a speech-controlled wireless communication device MS according to the flow chart of Fig. 2. In the speech recognition, an acoustic signal (speech) is converted, in a way known as such, into an electrical signal by a microphone, such as

15 a microphone 1a in the wireless communication device MS or a microphone 1b in a hands-free facility 2. The frequency response of the speech signal is typically limited to the frequency range below 10 kHz, *e.g.* in the frequency range from 100 Hz to 10 kHz. However, the frequency response of speech is not constant in the whole frequency

20 range, but there are more lower frequencies than higher frequencies. Furthermore, the frequency response of speech is different for different persons. In the method of the invention, the frequency range to be examined is divided into narrower sub-frequency ranges (M number of sub-bands). This is represented by block 101 in the appended Fig. 1.

25 These sub-frequency ranges are not made equal in width but taking into account the characteristic features of the speech, wherein some of the sub-frequency ranges are narrower and some are wider. At the low frequencies characteristic of speech, the division is denser, *i.e.* the sub-frequency ranges are narrower than for the higher frequencies, which

30 frequencies are more rare in speech. This idea is also applied in the Mel frequency scale, known as such, in which the width of frequency bands is based on the logarithmic function of frequency.

35 In connection with the division into sub-bands, the signals of the sub-bands are converted to a smaller sample frequency, *e.g.* by under-sampling or by low-pass filtering. Thus, samples are transferred from the block 101 to further processing at this lower sampling frequency. This sampling frequency is advantageously *ca.* 100 Hz, but it is obvious

that also other sampling frequencies can be applied within the scope of the present invention. These samples are converted into said feature vectors.

5 A signal formed in the microphone 1a, 1b is amplified in an amplifier 3a, 3b and converted into digital form in an analog-to-digital converter 4. The precision of the analog-to-digital conversion is typically in the range from 12 to 32 bits, and in the conversion of a speech signal, samples are taken advantageously 8'000 to 14'000 times a second, but the
10 invention can also be applied at other sampling rates. In the wireless communication device MS of Fig. 2, the sampling is arranged to be controlled by a controller 5. The audio signal in digital form is transferred to a speech recognition device 16 which is in a functional connection with the wireless communication device 16 and in which different
15 stages of the method according to the invention are processed. The transfer takes place *e.g.* via interface blocks 6a, 6b and an interface bus 7. In practical solutions the speech recognition device 16 can as well be arranged in the wireless communication device 16 itself or in another speech-controlled device, or as a separate auxiliary device or
20 the like.

The division into sub-bands is made preferably in a first filter block 8, to which the signal converted into digital form is conveyed. This first filter block 8 consists of several band-pass filters which are in this advantageous
25 embodiment implemented with digital technique and whose frequency ranges and band widths of the pass band differ from each other. Thus each band filtered part of the original signal passes the respective band-pass filter. For clarity, these band-pass filters are not shown separately in Fig. 2. These band-pass filters are implemented
30 advantageously in the application software of a digital signal processor (DSP) 13, which is known as such.

At the next stage 102, the number of sub-bands is reduced preferably by decimating in a decimating block 9, wherein L number of sub-bands
35 are formed ($L < M$), their energy levels being measurable. On the basis of the signal power levels of these sub-frequency ranges, it is possible to determine the signal energy in each sub-band. Also, the decimating

block 9 can be implemented in the application software of the digital signal processor 13.

5 An advantage obtained by the division into M sub-bands according to the block 1 is that the values of these M different sub-bands can be utilized in the recognition to verify the recognition result particularly in an application using coefficients according to the Mel frequency scale. However, the block 101 can also be implemented by forming directly L sub-bands, wherein the block 102 will not be necessary.

10 A second filter block 10 is provided for low pass filtering of signals of the sub-bands formed at the decimating stage (stage 103 in Fig. 1), wherein short changes in the signal strength are filtered off and they cannot have a significant effect in the determination of the energy level of the signal in further processing. After the filtration, a logarithmic function of the energy level of each sub-band is calculated in block 11 (stage 104) and the calculation results are stored for further processing in sub-band specific buffers formed in memory means 14 (not shown). These buffers are advantageously of the so-called FIFO type (First In -
15 First Out), in which the calculation results are stored as figures of *e.g.* 8 or 16 bits. Each buffer accommodates N calculation results. The value N depends on the application in question. Thus, the calculation results $p(t)$ stored in the buffer represent the filtered, logarithmic energy level of the sub-band at different measuring instants.

25 An arrangement block 12 performs so-called rank order filtering for the calculation results (stage 105), in which the mutual rank of the different calculation results are compared. At this stage 105, it is examined in the sub-bands whether there is possibly a pause in the speech. This examination is shown in a state machine chart in Fig. 3. The operations of this state machine are implemented substantially in the same way for each sub-band. The different functional states S0, S1, S2, S3 and S4 of the state machine are illustrated with circles. Inside these state circles are marked the operations to be performed in each functional state. The arrows 301, 302, 303, 304 and 305 illustrate the transitions from one functional state to another. In connection with these arrows are marked the criteria, whose realization will set off this transition. The curves 306, 307 and 308 illustrate the situation in which the functional
30
35

state is not changed. Also these curves are provided with the criteria for maintaining the functional state.

5 In the functional states S1, S2 and S3, a function $f()$ is shown, which represents the performing of the following operations in said functional states: preferably N calculation results $p(t)$ are stored in the buffer, and the lowest maximum value $p_min(t)$ and the highest minimum value $p_min(t)$ are determined advantageously by the following formulae:

$$10 \quad \begin{aligned} p_min(t) &= \min \left[\max \langle p(i - N + 1), p(i - N + 2), \dots, p(i) \rangle \right], & i = N, N + 1, \dots, t \\ p_max(t) &= \max \left[\min \langle p(i - N + 1), p(i - N + 2), \dots, p(i) \rangle \right], & i = N, N + 1, \dots, t \end{aligned}$$

Consequently, in the function $f(t)$, the maximum value $p_max(t)$ searched is the highest minimum value and the minimum value $p_min(t)$ is the lowest maximum value of the calculation results $p(i)$ stored in the different sub-band buffers. After this, the median power $p(t)_m$ is calculated, which is the median value of the calculation results $p(t)$ stored in the buffer, and a threshold value thr by the formula $thr = p_min + k \cdot (p_max - p_min)$, in which $0 < k < 1$. Next, in the function $f()$, a comparison is made between the median power $p(t)_m$ and the threshold value calculated above. The result of the calculation will set off different operations depending on the functional state in which the state machine is at a given time. This will be described in more detail hereinbelow in connection with the description of the different functional states.

After storing a group of sub-band specific calculation results $p(t)$ of the speech (N results per sub-band), the speech recognition device will move on to execute said state machine, which is implemented in the application software of either the digital signal processor 13 or the controller 5. The timing can be made in a way known as such, preferably with an oscillator, such as a crystal oscillator (not shown). The executing is started from the state $S0$, in which the variables to be used in the state machine are set in their initial values ($init()$): a pause counter C is set to zero, the power minimum p_min at the starting moment $t = 1$ ($p_min(t = 1)$) is set to the theoretical value of ∞ , in practice to the highest possible numerical value available in the speech recognition device.

This maximum value is influenced by the number of bits these power values are calculated with. Correspondingly, the power maximum p_{\max} at the starting moment $t = 1$ ($p_{\max}(t = 1)$) is set to the theoretical value of $-\infty$, in practice to the lowest possible numerical value
 5 available in the speech recognition device.

After setting of the initial values, the function moves on to the state S1, in which the operations of said function $f()$ are performed, wherein e.g. the power minimum p_{\min} and the power maximum p_{\max} as well as
 10 the median power $p(t)_m$ are calculated. In the functional state S1, also the pause counter C is increased by one. This functional state prevails until the expiry of a predetermined initial delay. This is determined by comparing the pause counter C with a predetermined beginning value BEG. At the stage when the pause counter C has reached the begin-
 15 ning value BEG, the operation moves on to state S2.

In the functional state S2, the pause counter C is set to zero and the operations of the function $f()$ are performed, such as storing of the new calculation result $p(t)$, and calculation of the power minimum p_{\min} , the
 20 power maximum p_{\max} as well as the median power $p(t)_m$ and the threshold value thr . The calculated threshold value and the median power are compared with each other, and if the median power is smaller than the threshold value, the operation moves on to state S3; in other cases, the functional state is not changed but the above-pre-
 25 sented operations of this functional state S2 are performed again.

In the functional state S3, the pause counter C is increased by one and the function $f()$ is performed. If the calculation indicates that the median power is still smaller than the threshold value, the value of the pause
 30 counter C is examined to find out if the median power has been below the power threshold value for a certain time. Expiry of this time limit can be found out by comparing the value of the pause counter C with an utterance time limit END. If the value of the counter is greater than or equal to said expiry time limit END, this means that no speech can be
 35 detected on said sub-band, wherein the state machine is exited.

However, if the comparison of the threshold value and the median power in the functional state S3 showed that the median power ex-

ceeded the power threshold value, it can be deduced that speech is detected on this sub-band, and the state machine returns to the functional state S2, in which *e.g.* the pause counter C is reset and the calculation is started from the beginning.

5

Consequently, the operation of a state machine to be used in the method according to an advantageous embodiment of the invention was described above in a general manner. In a speech recognition device according to the invention, the above-presented functional stages are performed separately for each sub-band.

10

Sampling a speech signal is performed advantageously at intervals, wherein the stages 101—104 are performed after the calculation of each feature vector, preferably at intervals of *ca.* 10 ms. Correspondingly, in the state machine of each sub-band, the operations according to the each active functional state are performed once (one calculation time), *e.g.* in state S3 the pause counter C(s) of the sub-band in question is increased, the function f(s) is performed, wherein *e.g.* a comparison is made between the median power and the threshold value, and on the basis of the same, the functional state is either retained or changed.

15

20

After one calculating round has been performed for the state machines of all the sub-bands, the operation moves on to stage 106 in the speech recognition, wherein it is examined on the basis of the information received from the different sub-bands whether a sufficiently long pause has been detected in the speech. This stage 106 is illustrated as a flow chart in the appended Fig. 4. For clarifying the examination, some comparison values are determined, which are given initial values preferably in connection with the manufacture of the speech recognition device, but if necessary, these initial values can be changed according to the application in question and the usage conditions. The setting of these initial values is illustrated with block 401 in the flow chart of Fig. 4:

25

30

35

- activity threshold SB_ACTIVE_TH whose value is greater than zero but smaller than the detection time limit END,
- detection quantity SB_SUFF_TH whose value is greater than zero but smaller than or equal to the number L of sub-bands,

- minimum number SB_MIN_TH of sub-bands whose value is greater than zero but smaller than the detection quantity SB_SUFF_TH.

5 In the method according to the invention, to detect a pause in speech it is examined, on how many sub-bands the energy level has possibly remained below said power threshold value and for how long. As disclosed in the functional description of the state machine above, the pause counter C indicates how long the audio energy level has re-
10 mained below the power threshold value. Thus, the value of the counter is examined for each sub-band. If the value of the counter is greater than or equal to the detection time limit END (block 402), this means that the energy level of the sub-band has remained below the power threshold value so long that a decision on detecting a pause can be
15 made for this sub-band, *i.e.* a sub-band specific detection is made. Thus, the detection counter SB_DET_NO is preferably increased by one.

20 If the value of the counter is greater than or equal to the activity threshold SB_ACTIVE_TH (block 404), the energy level on this sub-band has been below the power threshold value thr for a moment but not yet a time corresponding to the detection time limit END. Thus, the activity counter SB_ACT_NO in block 405 is increased preferably by one. In other cases, there is either an audio signal on the sub-band, or
25 the level of the audio signal has been below the power threshold value thr for only a short time.

Next, the operation moves on to block 406, in which the sub-band counter i used as an auxiliary variable is increased by one. On the ba-
30 sis of the value of this sub-band counter i, it can be deduced if all the sub-bands have been examined (block 407).

When the comparisons to the said pause counters have been made, it is examined, on how many sub-bands a pause was detected (the pause
35 counter was greater than or equal to the detection time limit END). If the number of such sub-bands is greater than or equal to the detection quantity SB_SUFF_TH (block 408), it is deduced in the method that there is a pause in the speech (pause detection decision, block 409),

and it is possible to move on to the actual speech recognition to find out what the user uttered. However, if the number of sub-bands is smaller than the detection quantity SB_SUFF_TH, it is examined, if the number of sub-bands including a pause is greater than or equal to the minimum number of sub-bands SB_MIN_TH (block 410). Furthermore, it is examined in block 411 if any of the sub-bands is active (the pause counter was greater than or equal to the activity threshold SB_ACTIVE_TH but smaller than the detection time limit END). In the method according to the invention, a decision is made in this situation that there is a pause in the speech if none of the sub-bands is active.

In a noise situation, noise on some sub-bands may effect that a detection decision cannot be made on all sub-bands even though there were a pause in the speech that should be detected. Thus, by means of said sub-band minimum SB_MIN_TH, it is possible to verify the detection of a pause in the speech particularly under noise conditions. Thus, in a noise situation, if a pause is detected on at least said minimum number SB_MIN_TH of sub-bands, a pause is detected in the speech if the pause detection decision on these sub-bands remains in force for the duration of said detection time limit END.

Correspondingly, under good conditions, using said detection time limit END may prevent a too quick decision on detecting a pause. Under good conditions, the said minimum number of sub-bands can quickly cause a pause detection decision, even though there is no such pause in the speech to be detected. By waiting the detection time limit for substantially all of the sub-bands, it is verified that there is actually a pause in the speech.

In another advantageous embodiment of the invention, it is not examined before making the decision of detecting a pause whether any of the sub-bands is active. Thus, the decision on detecting a pause is made on the basis of the results of the comparisons presented above.

The operations presented above can be implemented advantageously *e.g.* in the application software of the controller or digital signal processor of the speech recognition device.

The above-presented method for detecting a pause in speech according to the advantageous embodiment of the invention can be applied at the stage of teaching a speech recognition device as well as at the stage of speech recognition. At the teaching stage, the disturbance conditions can be usually kept relatively constant. However, when a speech-controlled device is used, the quantity of background noise and other interference can vary to a great extent. For improving the reliability of speech recognition particularly under varying conditions, the method according to another advantageous embodiment of the invention is supplemented with adaptivity to the calculation of the threshold value thr . For achieving this adaptivity, a modification coefficient $UPDATE_C$ is used, whose value is preferably greater than zero and smaller than one. The modification coefficient is first given an initial value within said value range. This modification coefficient is updated during speech recognition preferably in the following way. On the basis of the samples of the sub-bands stored in the buffers, a maximum power level win_max and a minimum power level win_min are calculated. After this, said calculated maximum power level win_max is compared with the power maximum p_max at the time, and said calculated minimum power level win_min is compared with the power minimum p_min . If the absolute value of the difference between the calculated maximum power level win_max and the power maximum p_max , or the absolute value of the difference between the calculated minimum power level win_min and the power minimum p_min has increased from the previous calculation time, the modification coefficient $UPDATE_C$ is increased. On the other hand, if the absolute value of the difference between the calculated maximum power level win_max and the power maximum p_max , or the absolute value of the difference between the calculated minimum power level win_min and the power minimum p_min has decreased from the previous calculation time, the modification coefficient $UPDATE_C$ is reduced. After this, a new power maximum and a new power minimum are calculated as follows:

$$\begin{aligned}
 p_min(t) &= (1 - UPDATE_C) \cdot p_min(t-1) + (UPDATE_C \cdot win_min) \\
 p_max(t) &= (1 - UPDATE_C) \cdot p_max(t-1) + (UPDATE_C \cdot win_max)
 \end{aligned}$$

The calculated new power maximum and minimum values are used at the next sampling round *e.g.* in connection with the performing of the function $f()$. The determination of this adaptive coefficient has *e.g.* the advantage that changes in the environmental conditions can be better taken into account in the speech recognition and the detection of a pause becomes more reliable.

The above-presented different operations for detecting a pause in the speech can be largely implemented in the application software of the controller and/or the digital signal processor of the speech recognition device. In the speech recognition device according to the invention, some of the functions, such as the division into sub-bands, can also be implemented with analog technique, which is known as such. In connection with the execution of the method, in the storing of the calculation results to be made at different stages, the variables, *etc.*, it is possible to use the memory means 14 of the speech recognition device, preferably a random access memory (RAM), a non-volatile random access memory (NVRAM), a FLASH memory, *etc.* The memory means 22 of the wireless communication device can as well be used for storing information.

Fig. 2, showing a the wireless communication device MS according to an advantageous embodiment of the invention, additionally shows a keypad 17, a display 18, a digital-to-analog converter 19, a headphone amplifier 20a, a headphone 21, a headphone amplifier 20b for a hands-free function 2, a headphone 21b, and a high-frequency block 23, all known *per se*.

The present invention can be applied in connection with several speech recognition systems functioning by different principles. The invention improves the reliability of detection of pauses in speech, which ensures the recognition reliability of the actual speech recognition. Using the method according to the invention, it is not necessary to perform the speech recognition in connection with a fixed time window, wherein the recognition delay is substantially not dependent on the rate at which the user utters speech commands. Also, the effect of background noise on speech recognition can be made smaller upon applying the method of the invention than is possible in speech recognition devices of prior art.

It is obvious that the invention is not limited solely to the embodiments presented above, but it can be modified within the scope of the appended claims.

Claims:

1. A method for detecting pauses in speech in speech recognition, in which method, for recognizing speech commands uttered by the user,
5 the voice is converted into an electrical signal, **characterized** in that in the method, the frequency spectrum of the electrical signal is divided into two or more sub-bands, samples of the signals in the sub-bands are stored at intervals, the energy levels of the sub-bands are determined on the basis of the stored samples, a power threshold value (thr)
10 is determined, and the energy levels of the sub-bands are compared with said power threshold value (thr), wherein the comparison results are used for producing a pause detecting result.

2. The method according to claim 1, **characterized** in that a detection
15 time limit (END) and a detection quantity (SB_SUFF_TH) are determined, wherein in the method, the calculation of the length of a pause in a sub-band is started when the energy level of the sub-band falls below said power threshold value (thr), wherein in the method, a sub-band specific detection is performed when the calculation reaches
20 the detection time limit (END), it is examined on how many sub-bands the energy level was below the power threshold value (thr) longer than the time detection limit (END), wherein a pause detection decision is made if the number of sub-band specific detections is greater than or equal to the detection quantity (SB_SUFF_TH).

- 25 3. The method according to claim 2, **characterized** in that in the method, also an activity time limit (SB_ACTIVE_TH) and an activity quantity (SB_MIN_TH) are determined, wherein a pause detection decision is made if the quantity of sub-band specific detections is greater
30 than or equal to the activity quantity (SB_MIN_TH) and the activity time limit (SB_ACTIVE_TH) has not been reached on the other sub-bands in the calculation of the length of the pause in the sub-band.

4. The method according to claim 1, 2 or 3, **characterized** in that the
35 power threshold value (thr) is calculated by the formula

$$thr = p_min + k \cdot (p_max - p_min), \text{ in which}$$

p_min = the smallest power maximum determined of the stored samples of the sub-bands, and
 p_max = the greatest power minimum determined of the stored samples of the sub-bands.

5

5. The method according to any of the claims 1 to 4, **characterized** in that said power threshold value (thr) is calculated adaptively by taking into account the environmental noise level at each instant.

10

6. The method according to claim 5, **characterized** in that for calculating said power threshold value (thr), a modification coefficient (UPDATE_C) is determined, and on the basis of the stored samples, the greatest power level (win_max) and the smallest power level (win_min) of the sub-bands are calculated, wherein the power maximum (p_max) and power minimum (p_min) are determined by the formulae:

15

$$p_max(i,t) = (1 - UPDATE_C) \cdot p_max(i,t-1) + (UPDATE_C \cdot win_max)$$

$$p_min(i,t) = (1 - UPDATE_C) \cdot p_min(i,t-1) + (UPDATE_C \cdot win_min)$$

20

in which $0 < UPDATE_C < 1$,
 $0 < i < L$, and
 L is the number of sub-bands.

25

7. The method according to claim 6, **characterized** in that further in the method,

30

— the modification coefficient (UPDATE_C) is increased, if the absolute value of the difference between said calculated highest power level (win_max) and the power maximum (p_max), or the absolute value of the difference between said calculated lowest power level (win_min) and the power minimum (p_min) has increased,

35

— the modification coefficient (UPDATE_C) is reduced, if the absolute value of the difference between said calculated highest power level (win_max) and the power maximum (p_max), or the absolute value of the difference between said calculated lowest power level (win_min) and the power minimum (p_min) has decreased.

8. A speech recognition device (16) comprising means (1a, 1b) for converting speech commands uttered by a user into an electrical signal, **characterized** in that it also comprises:

- 5 — means (8) for dividing the frequency spectrum of the electrical signal into two or more sub-bands,
- means (14) for storing samples of the signals of the sub-bands at intervals,
- means (5, 13) for determining energy levels of the sub-bands on the basis of the stored samples,
- 10 — means (5, 13) for determining a power threshold value (thr),
- means (5, 13) for comparing the energy levels of the sub-bands with said power threshold value (thr), and
- means (5, 13) for detecting a pause in the speech on the basis of said comparison results.
- 15

9. The speech recognition device (16) according to claim 8, **characterized** in that the power threshold value is calculated by the formula

20 $thr = p_min + k \cdot (p_max - p_min)$, in which

p_min = the smallest determined power maximum of the stored samples of the sub-bands, and

25 p_max = the greatest determined power minimum of the stored samples of the sub-bands.

10. The speech recognition device (16) according to claim 8 or 9, **characterized** in that it comprises also means (10, 11) for filtering the signals of the sub-bands before storage.

30

11. A wireless communication device (MS) comprising means (16) for recognizing speech and means (1a, 1b) for converting speech commands uttered by a user into an electrical signal, **characterized** in that the means (16) for recognizing speech comprise also:

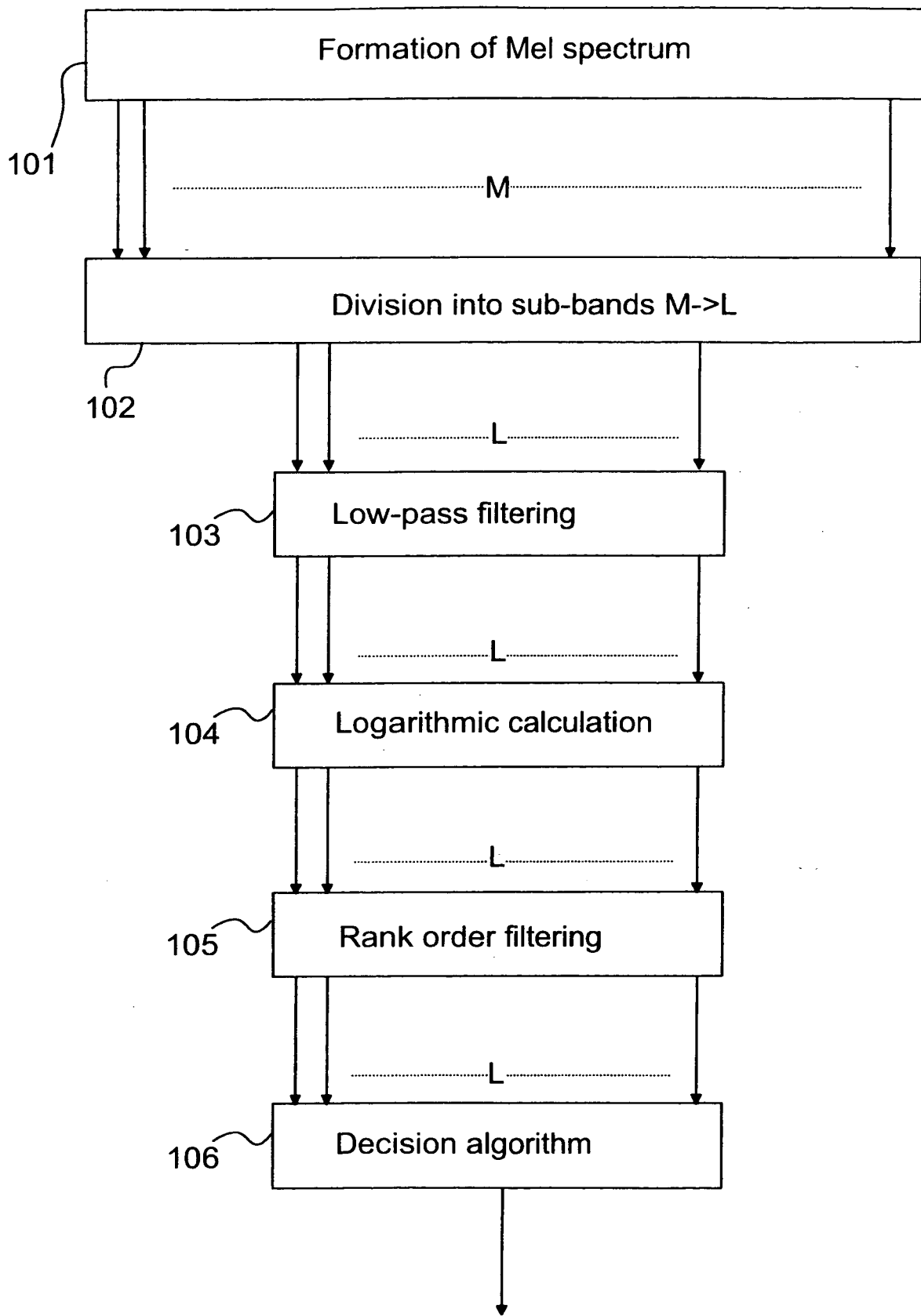
- 35 — means (8) for dividing the frequency spectrum of the electrical signal into two or more sub-bands,
- means (14) for storing samples of the signals of the sub-bands at intervals,

- means (5, 13) for determining energy levels of the sub-bands on the basis of the stored samples,
 - means (5, 13) for determining a power threshold value (thr),
 - means (5, 13) for comparing the energy levels of the sub-bands with said power threshold value (thr), and
 - means (5, 13) for detecting a pause in the speech on the basis of said comparison results.
- 5

Abstract

In a method for detecting pauses in speech in speech recognition, for recognizing speech commands uttered by the user, the voice is converted into an electrical signal, whose frequency spectrum is divided into two or more sub-bands. Samples of the signals on the sub-bands are stored at intervals, the energy levels of the sub-bands are determined on the basis of the stored samples, a power threshold value (thr) is determined, and the energy levels of the sub-bands are compared with said power threshold value (thr). The comparison results are used for producing a pause detecting result.

Fig. 1



JC584 U.S. PTO
09/482277
01/13/00

Fig 1

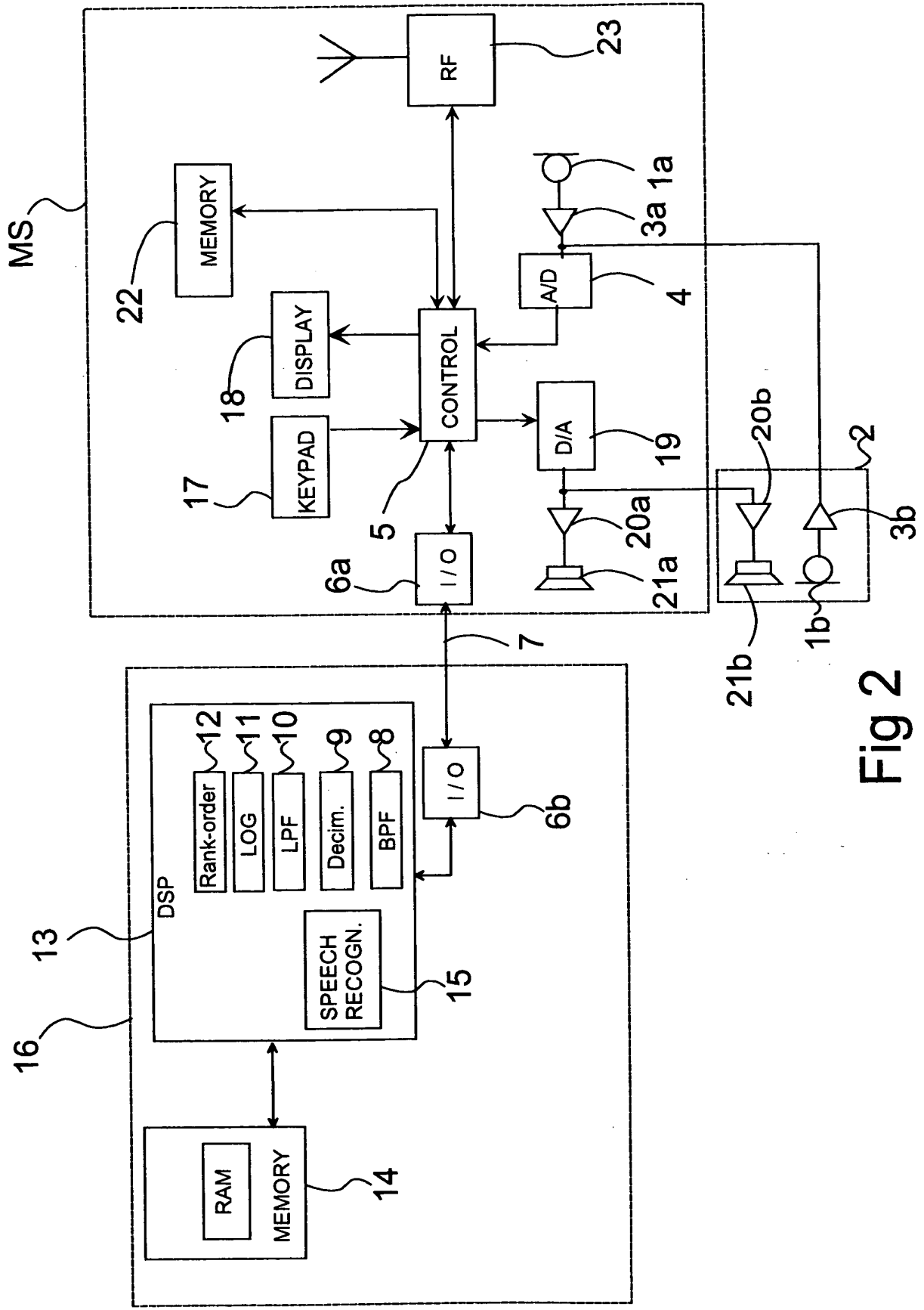


Fig 2

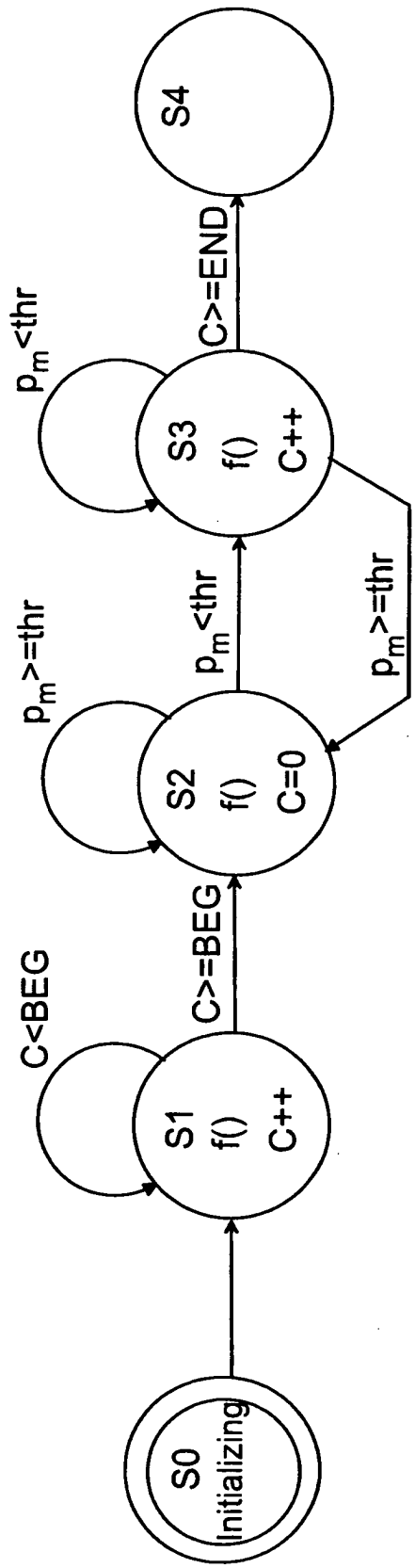


Fig 3

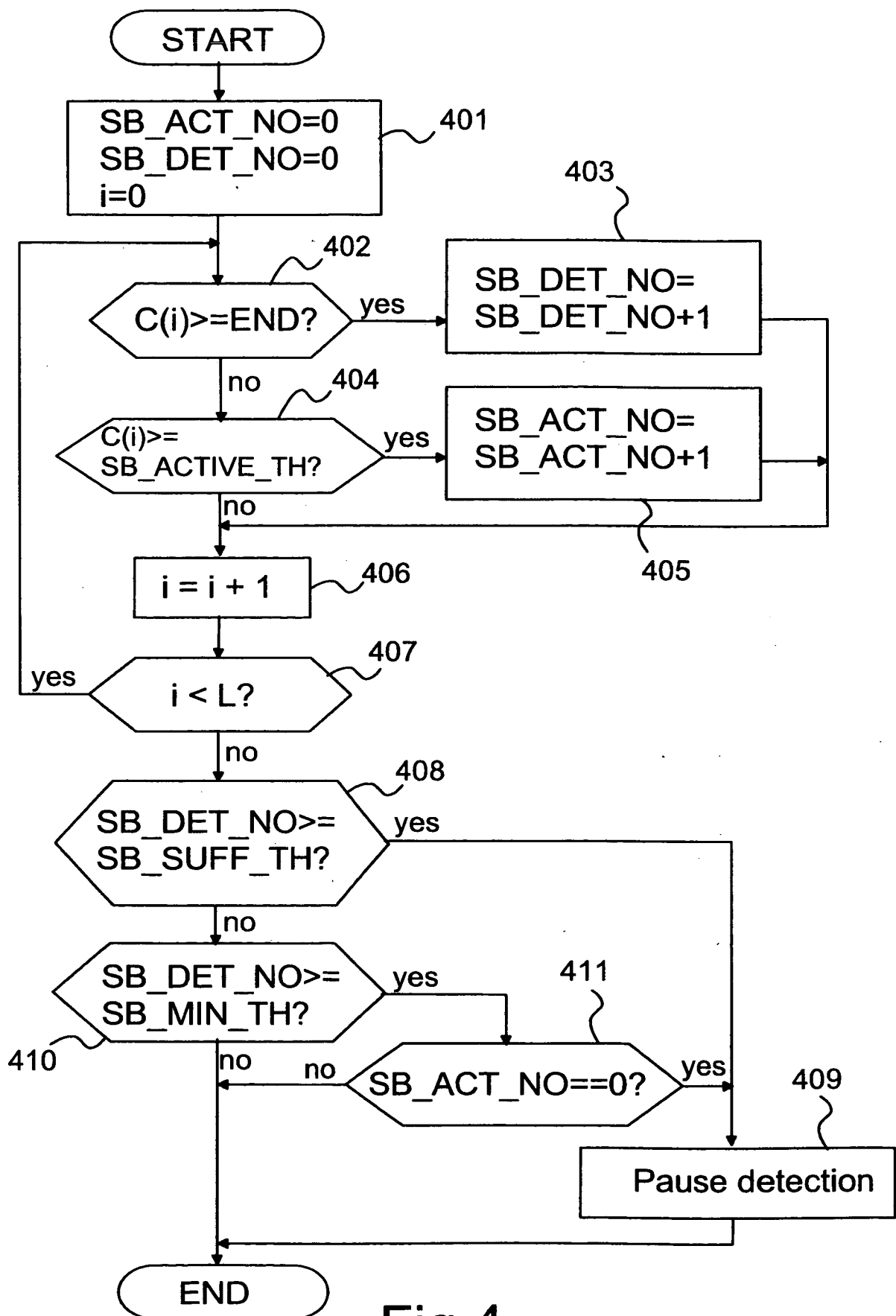


Fig 4

PATENTTI- JA REKISTERIHALLITUS
NATIONAL BOARD OF PATENTS AND REGISTRATION

Helsinki 19.11.1999

ETUOIKEUSTODISTUS
PRIORITY DOCUMENT

Hakija
Applicant

Nokia Mobile Phones Ltd
Espoo

Patenttihakemus nro
Patent application no

990078

Tekemispäivä
Filing date

18.01.1999

Kansainvälinen luokka
International class

G10L

Keksinnön nimitys
Title of invention

"Menetelmä puheen tunnistuksessa ja puheentunnistuslaite"

Täten todistetaan, että oheiset asiakirjat ovat tarkkoja jäljennöksiä patentti- ja rekisterihallitukselle alkuaan annetuista selityksestä, patenttivaatimuksista, tiivistelmästä ja piirustuksista.

This is to certify that the annexed documents are true copies of the description, claims, abstract and drawings originally filed with the Finnish Patent Office.


Pirjo Kaiffa
Tutkimussihteeri

Maksu 300,- mk
Fee 300,- FIM

Osoite: Arkadiankatu 6 A Puhelin: 09 6939 500 Telefax: 09 6939 5204
P.O.Box 1160 Telephone: + 358 9 6939 500 Telefax: + 358 9 6939 5204
FIN-00101 Helsinki, FINLAND

#2
JCS84 U.S. PTO
09/482277
01/13/00

Menetelmä puheentunnistuksessa ja puheentunnistuslaite

5 Nyt esillä oleva keksintö kohdistuu oheisen patenttivaatimuksen 1 johdanto-osan mukaiseen menetelmään puheentunnistuksessa, oheisen patenttivaatimuksen 8 johdanto-osan mukaiseen puheentunnistuslaitteeseen ja oheisen patenttivaatimuksen 11 johdanto-osan mukaiseen puheella ohjattavaan langattomaan viestimeen.

10 Langattomien viestimien käytön helpottamiseksi on kehitetty puheentunnistuslaitteita, joiden avulla käyttäjä voi lausua puhekomentoja, jotka puheentunnistuslaite pyrkii tunnistamaan ja muuntamaan puhekomentoa vastaavaksi toiminnoksi, esim. puhelinnumeron valintakomennoksi. Hankaluutena puheohjauksen toteuttamisessa on mm. se, että eri
15 käyttäjät lausuvat puhekomennot eri tavalla: puhenopeus voi olla erilainen eri käyttäjillä, samoin puheen voimakkuus, äänen sävy jne. Lisäksi puheentunnistusta häiritsee mahdollinen taustamelu, jonka häiritsevyys ulkona ja autossa voi olla huomattavaa. Taustamelu vaikeuttaa sanojen tunnistusta sekä eri sanojen erottamista toisistaan esim. puhelinnumeroa lausuttaessa.
20

Joissakin puheentunnistuslaitteissa on käytetty kiinteään aika-ikkunaan perustuvaa tunnistusmenetelmää. Tällöin käyttäjällä on ennalta määrätty aika, jonka kuluessa hänen on lausuttava haluamansa komentosana.
25 Aika-ikkunan kuluttua umpeen puheentunnistuslaite pyrkii selvittämään, minkä sanan/komennon käyttäjä lausui. Tällaiseen kiinteään aika-ikkunaan perustuvassa menetelmässä on kuitenkin mm. se epäkohta, että kaikki lausuttavat sanat eivät ole yhtä pitkiä, esim. nimien kohdalla etunimi on usein selvästi lyhyempi kuin sukunimi. Tällöin lyhyemmän sanan jälkeen kuluu enemmän aikaa tunnistukseen kuin pidemmän sanan tunnistuksessa. Tämä on epämiellyttävää käyttäjän kannalta. Lisäksi aika-ikkuna on asetettava hitaampien puhujien mukaan, ettei tunnistusta aloiteta, ennen kuin koko sana on lausuttu. Nopeammin sanoja lausuttaessa viive lausumisen ja tunnistuksen välillä
30 lisää epämiellyttävyyden tunnetta.
35

Toinen tunnettu puheentunnistusmenetelmä perustuu puhesignaaleista muodostettuihin malleihin ja niiden vertailuun. Kommentosanoista muo-

dostetut mallit on etukäteen tallennettu tai käyttäjä on voinut opettaa haluamiaan sanoja, joista on muodostettu ja tallennettu mallit. Puheentunnistuslaite vertailee tallennettuja malleja käyttäjän lausumista ään-teistä muodostettuihin piirvektoreihin sanojen lausumisen aikana ja
5 laskee todennäköisyyksiä puheentunnistuslaitteen sanaston eri sanoille (komentosanoille). Todennäköisyyden ylittäessä jollakin komentosanal-la ennalta asetetun arvon, puheentunnistuslaite valitsee tämän komen-tosanan tunnistustulokseksi. Tällöin voi virheellisiä tunnistustuloksia syntyä erityisesti sellaisten sanojen kohdalla, joissa sanan alku muistut-
10 taa äänteellisesti jotakin muuta sanastoon kuuluvaa sanaa. Esimerkiksi käyttäjä on opettanut puheentunnistuslaitteelle sanat "Mari" ja "Marika". Jos käyttäjä lausuu sanaa "Marika", saattaa puheentunnistuslaite tehdä tunnistuspäätökseksi "Mari", vaikka käyttäjä ei olisi ehtinyt lausua vielä sanan loppua. Tällaisissa puheentunnistuslaitteissa käytetään usein ns.
15 Hidden-Markov-Model -puheentunnistusmenetelmää (HMM).

Patentissa US-4,870,686 on esitetty puheentunnistusmenetelmä ja pu-heentunnistuslaite, jossa käyttäjän sanojen lopun ilmaiseminen perus-tuu hiljaisuuteen, siis puheentunnistuslaite tutkii, onko äänisignaalia
20 havaittavissa vai ei. Ongelmana tässä ratkaisussa on se, että liian voi-makas taustamelu voi estää taukojen havaitsemisen, jolloin puheen-tunnistus ei onnistu.

Nyt esillä olevan keksinnön eräänä tarkoituksena on aikaansaada pa-rannettu menetelmä puheessa olevien taukojen havaitsemiseksi ja pu-heentunnistuslaite. Keksintö perustuu siihen ajatukseen, että jaetaan
25 tutkittava äänikaista alikaistoihin ja tutkitaan signaalin tehoa kullakin alikaistalla. Mikäli riittävän usealla alikaistalla signaalin teho alittaa tie-tyn rajan riittävän pitkän ajan, tehdään päätelmä siitä, että puheessa on tauko. Nyt esillä olevan keksinnön mukaiselle menetelmälle on tunnus-
30 omaista se, mitä on esitetty oheisen patenttivaatimuksen 1 tunnus-merkkiosassa. Nyt esillä olevan keksinnön mukaiselle puheentunnistus-laitteelle on tunnusomaista se, mitä on esitetty oheisen patenttivaati-muksen 8 tunnusmerkkiosassa. Nyt esillä olevan keksinnön mukaiselle
35 langattomalle viestimelle on tunnusomaista se, mitä on esitetty oheisen patenttivaatimuksen 11 tunnusmerkkiosassa.

Nyt esillä olevalla keksinnöllä saavutetaan merkittäviä etuja tunnetun tekniikan mukaisiin ratkaisuihin verrattuna. Keksinnön mukaisella menetelmällä saadaan luotettavampi sanavälin ilmaisu kuin tunnetun tekniikan mukaisilla menetelmillä. Tällöin puheentunnistuksen luotettavuus
5 paranee ja virheellisten tunnistusten ja epäonnistuneiden tunnistusten määrä pienenee. Lisäksi puheentunnistustaite on joustavampi erilaisten käyttäjien puhetottumusten suhteen, koska puhekomennot voidaan lausua hitaammin tai nopeammin ilman, että tunnistuksessa on epämiellyttävää viivettä tai että tunnistus tapahtuisi kesken sanan lausumisen.

10 Keksinnön mukaisella alikaistoihin jakamisella saadaan ulkoisten häiriöiden vaikutusta pienennettyä. Tyypillisesti häiriösignaalit esim. autossa ovat suhteellisen matalataajuisia. Tunnetun tekniikan mukaisissa ratkaisuissa koko käsiteltävän signaalin taajuusalueen sisältämää energiaa käytetään tunnistuksessa hyväksi, jolloin voimakkaat mutta kapeakaistaiset signaalit heikentävät signaali-kohinasuhdetta merkittävästi. Sen sijaan jaettaessa tutkittava taajuusalue keksinnön mukaisesti alikaistoihin, saadaan sellaisilla alikaistoilla, joilla häiritsevien signaalien osuus on suhteellisen pieni, signaali-kohinasuhdetta parannettua merkittävästi, mikä parantaa tunnistusvarmuutta.
20

Nyt esillä olevaa keksintöä selostetaan seuraavassa tarkemmin viitaten samalla ohjeisiin piirustuksiin, joissa

25 kuva 1 esittää vuokaaviona keksinnön erään edullisen suoritusmuodon mukaista menetelmää,

kuva 2 esittää keksinnön erään edullisen suoritusmuodon mukaista puheentunnistustaitea pelkistettynä lohkokaaaviona,
30

kuva 3 esittää keksinnön erään edullisen suoritusmuodon mukaisessa menetelmässä sovellettavaa sijalukusuodatusta (rank-order filtering) tilakonekaaviona, ja

35 kuva 4 esittää vuokaaviona keksinnön erään edullisen suoritusmuodon mukaisessa menetelmässä sovellettavaa tauon päättelylogiikkaa.

Selostetaan seuraavassa keksinnön erään edullisen suoritusmuodon mukaisen menetelmän toimintaa viitaten samalla kuvan 1 vuokaavioon käyttäen esimerkkinä kuvan 2 lohkoakaavion mukaista puheella ohjattavaa langatonta viestintä MS. Puheentunnistuksessa suoritetaan si-
5 nänsä tunnetusti akustisen signaalin (puheen) muuntaminen sähköiseksi signaaliksi mikrofonilla, kuten langattoman viestimen MS mikrofonilla 1a tai kaiutintoiminnon 2 mikrofonilla 1b. Puhesignaalin taajuusvaste rajoittuu tyypillisesti alle 10 kHz:n taajuusalueelle, esim. taajuus-
10 alueelle 100 Hz—10 kHz. Puheen taajuusvaste ei kuitenkaan ole vakio koko taajuusalueella, vaan siinä matalampia taajuuksia esiintyy enemmän kuin korkeampia taajuuksia. Lisäksi eri henkilöillä puheen taajuusvaste on erilainen. Keksinnön mukaisessa menetelmässä tutkittava taajuusalue jaetaan kapeampiin alitaajuusalueisiin (alikaistoihin, M kpl). Tätä esittää lohko 101 oheisessa kuvassa 1. Näitä alitaajuusalueita ei
15 tehdä tasalevyisiksi, vaan puheen ominaispiirteet huomioiden, jolloin osa alitaajuusalueista on kapeampia ja osa on leveämpiä. Puheelle ominaisilla, alemmilla taajuuksilla jako on tiheämpi, eli alitaajuusalueet ovat kapeampia, kuin puheessa harvemmin esiintyvillä, korkeammilla taajuuksilla. Tähän perustuu myös sinänsä tunnettu mel-taajuusjako
20 (Mel Frequency Scale), jossa taajuuskaistojen leveys perustuu logaritmiseen taajuuden funktioon.

Allikaistoihin jakamisen yhteydessä alikaistojen signaalit muunnetaan pienemmälle näytetaajuudelle esim. alinäytteistämällä tai alipäästösuo-
25 dattamalla. Tällöin lohkoista 101 näytteitä siirretään jatkokäsittelyyn tällä alemmalla näytetaajuudella. Tämä näytetaajuus on edullisesti n. 100 Hz, mutta on selvää, että nyt esillä olevan keksinnön puitteissa myös muita näytetaajuuksia voidaan soveltaa. Näistä näytteistä muodostetaan mainittuja piirrevektoreita.

30 Mikrofonissa 1a, 1b muodostettu signaali vahvistetaan vahvistimessa 3a, 3b ja muunnetaan digitaaliseksi analogia-digitaalimuuntimessa 4. Analogia/digitaalimuunnoksen tarkkuus on tyypillisesti välillä 12—32 bittiä ja puhesignaalin muuntamisessa näytteitä otetaan edullisesti
35 8000—14000 kertaa sekunnissa, mutta keksintöä voidaan soveltaa myös muilla näytteenottonopeuksilla. Kuvan 2 langattomassa viestimessä MS näytteenotto on järjestetty suoritettavaksi kontrollerin 5 ohjaamana. Digitaalisessa muodossa oleva äänisignaali siirretään langat-

5 toman viestimen MS kanssa toiminnallisessa yhteydessä olevaan puheentunnistuslaitteeseen 16, jossa suoritetaan keksinnön edullisen suoritusbuodon mukaisen menetelmän eri vaiheita. Siirto suoritetaan esim. liityntälohkojen 6a, 6b ja liityntäväylän 7 kautta. Puheentunnistus-
5 laite 16 voi käytännön sovelluksissa olla toteutettuna myös itse langattomassa viestimessä MS tai muussa puheohjattavassa laitteessa, tai erillisenä lisälaitteena tai vastaavana.

10 Alikeistoihin jako tehdään edullisesti ensimmäisessä suodatinlohkossa 8, johon digitaaliseksi muunnettu signaali johdetaan. Tämä ensimmäinen suodatinlohko 8 koostuu useista, tässä edullisessa suoritusbuodossa digitaalitekniikalla toteutetuista, kaistanpäästösuodattimista, joiden päästökaistan taajuusalueet sekä kaistanleveydet eroavat toisistaan. Tällöin kunkin kaistanpäästösuodattimen läpäisee alkuperäisestä
15 signaalista kaistanpäästösuodatettu osa. Selvyyden vuoksi ei kuvassa 2 ole esitetty erillisinä näitä kaistanpäästösuodattimia. Nämä kaistanpäästösuodattimet on toteutettu edullisesti signaalinkäsittely-yksikön 13 (DSP, Digital Signal Processor) sovellusohjelmistossa, kuten on siinä tunnettua.

20 Seuraavassa vaiheessa 102 vähennetään alikaistojen lukumäärää edullisesti desimoimalla desimointilohkossa 9, jolloin muodostuu L kappaletta alikaistoja ($L < M$), joiden energiatasot ovat mitattavissa. Näiden alitaajuusalueiden signaalinvoimakkuuksien perusteella voidaan määrittää signaalin energia kullakin alikaistalla. Myös desimointilohko 9 voidaan toteuttaa digitaalisen signaalinkäsittely-yksikön 13 sovellusohjelmistossa.

30 Etu, joka saavutetaan lohkon 1 mukaisella M alikaistalla jakamisella on se, että näitä M:n eri alikaistan arvoja voidaan käyttää tunnistuksessa apuna tunnistustuloksen varmentamiseksi erityisesti sellaisessa sovelluksessa, jossa käytetään Mel-taajuusjaon mukaisia kertoimia. Lohko 101 voidaan kuitenkin toteuttaa myös siten, että siinä muodostetaan suoraan L kappaletta alikaistoja, jolloin lohkoa 102 ei tarvita.

35 Toisessa suodatinlohkossa 10 suoritetaan desimointivaiheessa muodostetuille alikaistojen signaaleille alipäästösuodatus (vaihe 103 kuvassa 1), jolloin lyhyet signaalinvoimakkuuden muutokset suodattuvat

ja eivät pääse vaikuttamaan merkittävästi signaalin energiatason määrittämiseen jatkossa. Suodatuksen jälkeen lasketaan lohossa 11 kunkin alikaistan energiatasosta logaritmifunktio (vaihe 104), jonka muodostamat laskentatulokset tallennetaan jatkokäsittelyä varten muistivä-

5 lineisiin 14 muodostettuihin alikaistakohtaisiin puskuireihin (ei esitetty). Nämä puskurit ovat edullisesti ns. FIFO-tyyppisiä (First In - First Out), joihin laskentatulokset tallennetaan esim. 8- tai 16-bittisinä lukuina. Kunkin puskuriin mahtuu N kappaletta laskentatuloksia. Arvo N riippuu kulloisestakin sovelluksesta. Puskuuriin tallennetut laskentatulokset $p(t)$

10 kuvaavat siis alikaistan suodatettua, logaritmista energiatasoa eri mitausajanhetkinä.

Järjestelylohko 12 suorittaa laskentatuloksille ns. rank-order-suodatuksen (vaihe 105), jossa eri laskentatulosten keskinäistä suuruutta vertail-

15 laan. Tässä vaiheessa 105 tutkitaan alikaistoittain se, onko puheessa mahdollisesti tauko. Tämä tutkiminen on esitetty tilakonekaaviona kuvassa 3. Tämän tilakoneen toiminnot toteutetaan olennaisesti samalaisina kullekin alikaistalle. Tilakoneen eri toimintatiloja S0, S1, S2, S3 ja S4 on esitetty ympyröillä. Näiden tilaympyröiden sisään on merkitty kussakin toimintatilassa suoritettavat toimenpiteet. Nuolet 301, 302, 20 303, 304 ja 305 kuvaavat siirtymisiä toimintatiloista toiseen. Näiden nuolien yhteyteen on merkitty kriteerit, joiden toteutuminen aikaansaa tämän siirtymisen. Kaaret 306, 307 ja 308 kuvaavat tilannetta, jossa toimintatilaa ei vaihdeta. Myös näiden kaarien yhteyteen on merkitty

25 kriteerit toimintatilan säilyttämiseksi ennallaan.

Toimintatiloissa S1, S2 ja S3 on esitetty funktio $f()$, joka tarkoittaa seuraavien toimenpiteiden suorittamista mainituissa toimintatiloissa: laskentatuloksia $p(t)$ tallennetaan puskuuriin edullisesti N kappaletta, joista 30 etsitään pienin maksimiarvo $p_{\min}(t)$ ja suurin minimiarvo $p_{\min}(t)$ edullisesti seuraavilla kaavoilla:

$$p_{\min}(t) = \min \left[\max \left(p(i - N + 1), p(i - N + 2), \dots, p(i) \right) \right], \quad i = N, N + 1, \dots, t$$

$$p_{\max}(t) = \max \left[\min \left(p(i - N + 1), p(i - N + 2), \dots, p(i) \right) \right], \quad i = N, N + 1, \dots, t$$

35

Funktiossa $f()$ haetaan siis maksimiarvoksi $p_{\max}(t)$ eri alikaistapuskuireihin tallennetuista laskentatuloksista $p(i)$ suurin minimiarvo ja mi-

- nimiärvoksi $p_{\min}(t)$ pienin maksimiärv. Tämän jälkeen lasketaan mediaaniteho $p(t)_m$, joka on mediaaniärv puskuriin tallennetuista laskentatuloksista $p(t)$ sekä kynnysärv thr kaavalla $thr = p_{\min} + k \cdot (p_{\max} - p_{\min})$, jossa $0 < k < 1$. Seuraavaksi funktiossa $f()$ suoritetaan mediaanitehon $p(t)_m$ vertailu edellä laskettuun kynnysärvöön. Vertailun tulos aikaansaa erilaisia toimenpiteitä riippuen siitä, missä toimintatilassa tilakone kulloinkin on. Tätä kuvataan jäljempänä tarkemmin eri toimintatilojen kuvauksen yhteydessä.
- 10 Sen jälkeen kun puheesta on tallennettu joukko alikaistakohtaisia laskentatuloksia $p(t)$ (N kpl/alikaista), puheentunnistuslaite siirtyy suorittamaan mainittua tilakoneetta, joka on toteutettu joko digitaalisen signaalinkäsittely-yksikön 13 tai kontrollerin 5 sovellusohjelmistossa. Ajoitus voidaan muodostaa sinänsä tunnetusti edullisesti oskillaattorilla, kuten
- 15 kideoskillaattorilla (ei esitetty). Suoritus aloitetaan tilasta S_0 , jossa tehdään tilakoneessa käytettävien muuttujien asettamiset alkuärvöihin ($init()$): taukolaskuri C nollataan, tehominimiärv p_{\min} aloitusajanhetkellä $t=1$ ($p_{\min}(t=1)$) asetetaan teoreettisesti ärvöön ∞ , käytännössä puheentunnistuslaitteessa käytettävissä olevaksi suurimmaksi mahdolliseksi lukuärvoksi. Tähän maksimiärvöön vaikuttaa se, kuinka monella bitillä näitä tehoärvöjä lasketaan. Vastaavasti tehomaksimiärv p_{\max} aloitusajanhetkellä $t=1$ ($p_{\max}(t=1)$) asetetaan teoreettisesti ärvöön $-\infty$, käytännössä puheentunnistuslaitteessa käytettävissä olevaksi pienimmäksi mahdolliseksi lukuärvoksi.
- 20
- 25 Alkuärvöjen asetuksen jälkeen toiminta siirtyy tilaan S_1 , jossa suoritetaan mainitun funktion $f()$ edellä esitetyt toimenpiteet, jolloin mm. tehojen minimiärv p_{\min} ja maksimiärv p_{\max} sekä mediaaniteho $p(t)_m$ lasketaan. Toimintatilassa S_1 kasvatetaan lisäksi taukolaskuria C yhdellä. Tässä toimintatilassa pysytään, kunnes ennalta määritetty alkuviive on kulunut umpeen. Tämä selvitetään vertailemalla taukolaskuria C ennalta asetettuun aloitusärvöön BEG . Siinä vaiheessa kun taukolaskuri C on saavuttanut aloitusärvön BEG , toiminta siirtyy tilaan S_2 .
- 30
- 35 Toimintatilassa S_2 taukolaskuri C nollataan ja suoritetaan funktion $f()$ toimenpiteet, kuten uuden laskentatuloksen $p(t)$ tallennus, tehominimin p_{\min} , tehomaksimin p_{\max} ja mediaanitehon $p(t)_m$ sekä kynnysärvön thr laskenta. Laskettua kynnysärvöä ja mediaanitehoä verrataan kes-

kenään ja mikäli mediaaniteho on pienempi kuin kynnysarvo, siirrytään toimintatilaan S3, muussa tapauksessa toimintatilaa ei vaihdeta, vaan suoritetaan tämän toimintatilan S2 edellä esitetyt toimenpiteet uudelleen.

5

Toimintatilassa S3 kasvatetaan taukolaskuria C yhdellä ja suoritetaan funktio f(). Jos vertailu osoittaa, että mediaaniteho on edelleen pienempi kuin kynnysarvo, tutkitaan taukolaskurin C arvo sen selvittämiseksi, onko mediaaniteho ollut tietyn ajan alle tehon kynnysarvon. Tämän aikarajan täytyminen on selvitettävissä vertaamalla taukolaskurin C arvoa ilmaisuaikarajaan END. Jos laskurin arvo on suurempi tai yhtäsuuri kuin mainittu ilmaisuaikaraja END, merkitsee se sitä, että kyseisellä alikaistalla ei puhetta ole havaittavissa, jolloin poistutaan tilakoneesta.

15

Jos toimintatilassa S3 kynnysarvon ja mediaanitehon vertailu kuitenkin osoitti, että mediaaniteho on ylittänyt tehon kynnysarvon, voidaan tästä tehdä päätelmä, että puhetta on tällä alikaistalla havaittavissa ja tilakone palautuu toimintatilaan S2, jossa mm. taukolaskuri C nollataan ja laskenta aloitetaan alusta.

20

Edellä oli siis kuvattu keksinnön erään edullisen suoritusmuodon mukaisessa menetelmässä käytettävän tilakoneen toimintaa yleisesti. Keksinnön mukaisessa puheentunnistusalitteessa edellä esitetyt toimintavaiheet suoritetaan kunkin alikaistan osalta erikseen.

25

Näytteenotto puhesignaalista suoritetaan edullisesti määräväleillä, jolloin vaiheet 101—104 suoritetaan kunkin piirvektorin laskennan jälkeen, edullisesti n. 10 ms:n välein. Vastaavasti kunkin alikaistan tilakoneessa suoritetaan kulloinkin aktiivisena olevan toimintatilan mukaiset toimenpiteet kerran (yksi laskentakierros), esim. tilassa S3 kasvatetaan ao. alikanavan taukolaskuria C(s), suoritetaan funktio f(s), jossa mm. tehdään mediaanitehon ja kynnysarvon välinen vertailu ja sen perusteella joko säilytetään toimintatila ennallaan tai muutetaan toimintatilaa.

30

Kun kaikkien alikaistojen tilakoneiden osalta on suoritettu yksi laskentakierros, siirrytään puheentunnistuksessa vaiheeseen 106, jossa tutkitaan eri alikaistoista saadun informaation perusteella se, onko puhees-

sa havaittu riittävän pitkä tauko. Tätä vaihetta 106 on kuvattu vuokaaviona oheisessa kuvassa 4. Tutkimisen selventämiseksi määritetään muutamia vertailuarvoja, joille annetaan alkuarvot edullisesti puheentunnistulaitteen valmistuksen yhteydessä, mutta näitä alkuarvoja voidaan tarvittaessa muuttaa kulloisenkin sovelluksen ja käyttöolosuhteiden mukaan. Näiden alkuarvojen asettamista esittää lohko 401 kuvan 4 vuokaaviossa:

- aktiivisuuskynnys SB_ACTIVE_TH, jonka arvo on suurempi kuin nolla, mutta pienempi kuin ilmaisuaikaraja END;
- 10 – ilmaisumäärä SB_SUFF_TH, jonka arvo on suurempi kuin nolla, mutta pienempi tai yhtäsuuri kuin alikaistojen lukumäärä L,
- alikaistojen minimimäärä SB_MIN_TH, jonka arvo on suurempi kuin nolla, mutta pienempi kuin ilmaisumäärä SB_SUFF_TH.

- 15 Keksinnön mukaisessa menetelmässä puheessa olevan tauon havaitsemiseksi tutkitaan, kuinka monella alikaistalla energiataso on mahdollisesti pysynyt mainitun tehon kynnysarvon alapuolella ja kuinka kauan. Kuten edellä olevasta tilakoneen toimintakuvauksesta käy ilmi, tauokaskuri C ilmaisee sen, kuinka pitkään alikaistalla on äänen energiataso
- 20 ollut tehon kynnysarvon alapuolella. Tällöin tutkitaan kunkin alikaistan laskurin arvoa. Jos laskurin arvo on suurempi tai yhtä suuri kuin ilmaisuaikaraja END (lohko 402), merkitsee se sitä, että alikaistan energiataso on ollut tehon kynnysarvon alapuolella niin kauan, että päätös tauon havaitsemisesta voidaan tehdä tämän alikaistan osalta, eli muodostetaan alikanavakohtainen ilmaisu. Tällöin lohkossa 403 kasvatetaan ilmaisulaskuria SB_DET_NO edullisesti yhdellä.
- 25

- Jos laskurin arvo on suurempi tai yhtä suuri kuin aktiivisuuskynnys SB_ACTIVE_TH (lohko 404), energiataso tällä alikaistalla on ollut tehon kynnysarvon thr alapuolella hetken, mutta ei vielä ilmaisuaikarajaa
- 30 END vastaavaa aikaa. Tällöin lohkossa 405 kasvatetaan aktiivisuuskaskuria SB_ACT_NO edullisesti yhdellä. Muussa tapauksessa alikaistassa on joko äänisignaalia, tai äänisignaalin taso on ollut vain lyhyen ajan alle tehon kynnysarvon thr.

35

Seuraavaksi siirrytään lohkoon 406, jossa apumuuttujana käytettävää alikaistalaskuria i kasvatetaan yhdellä. Tämän alikaistalaskurin i arvon

- Kun vertaillaan mainittuihin taukolaskureihin suoritettu, tutkitaan, kuinka monella alikaistalla on havaittu tauko (taukolaskuri oli suurempi tai yhtäsuuri kuin ilmaisuaikaraja END). Jos tällaisten alikaistojen lukumäärä on suurempi tai yhtäsuuri kuin ilmaisumäärä SB_SUFF_TH (lohko 408), menetelmässä päätellään, että puheessa on tauko (tauon tunnistuspäätös, lohko 409) ja voidaan siirtyä varsinaiseen puheentunnistukseen, jossa pyritään selvittämään se, mitä käyttäjä lausui. Jos sen sijaan alikaistojen lukumäärä on pienempi kuin ilmaisumäärä SB_SUFF_TH, tutkitaan, onko alikaistojen, joissa on tauko, määrä suurempi tai yhtäsuuri kuin alikaistojen minimimäärä SB_MIN_TH (lohko 410). Lohkossa 411 tutkitaan vielä, onko jokin alikaista aktiivinen (taukolaskuri oli suurempi tai yhtäsuuri kuin aktiivisuuskynnys SB_ACTIVE_TH, mutta pienempi kuin ilmaisuaikaraja END). Keksinnön mukaisessa menetelmässä tehdään tässä tilanteessa päätös siitä, että puheessa on tauko, jos mikään alikaista ei ole aktiivinen.
- 20 Kohinatilanteessa voi joillakin alikaistoilla kohina vaikuttaa siten, että ilmaisupäätöstä ei saada kaikilla alikaistoilla, vaikka puheessa olisi tauko, joka tulisi ilmaista. Tällöin mainitun alikaistojen minimimäärän SB_MIN_TH avulla voidaan puheessa olevan tauon ilmaisua varmentaa erityisesti kohinaisissa olosuhteissa. Tällöin kohinatilanteessa, mikäli tauko havaitaan vähintään mainitulla minimimäärällä SB_MIN_TH alikaistoja, todetaan puheessa oleva tauko, jos tauon havaitsemispäätös näillä alikaistoilla pysyy voimassa mainitun ilmaisuaikarajan END verran.
- 30 Vastaavasti hyvissä olosuhteissa mainitun ilmaisuaikarajan END käyttämisellä voidaan estää liian nopea tauon ilmaisupäätös. Hyvissä olosuhteissa voi mainitulla minimimäärällä alikaistoja tauon ilmaisupäätös tulla hyvinkin nopeasti, vaikka puheessa ei olisi sellaista taukoa, joka tulisi ilmaista. Odottamalla olennaisesti kaikkien alikanavien osalta ilmaisuaikarajan verran varmennetaan sitä, että puheessa todella on tauko.

Keksinnön eräässä toisessa edullisessa suoritusmuodossa ei ennen tauon tunnistuspäätöksen tekemistä tutkita sitä, onko jokin alikaista aktiivinen. Tällöin tauon tunnistuspäätös tehdään edellä esitettyjen vertailujen tuloksien perusteella.

5

Edellä esitetyt toiminnot voidaan edullisesti toteuttaa esimerkiksi puheentunnistulaitteen kontrollerin tai digitaalisen signaalinkäsittely-yksikön sovellusohjelmistossa.

- 10 Edellä esitettyä keksinnön edullisen suoritusmuodon mukaista menetelmää puheessa olevan tauon ilmaisemiseksi voidaan soveltaa puheentunnistulaitteen opetusvaiheessa sekä puheentunnistusvaiheessa. Opetusvaiheessa voidaan häiriöolosuhteet pitää tavallisesti suhteellisen vakioina. Sen sijaan käytettäessä puheella ohjattavaa laitetta voi taustamelun ja muiden häiriöiden määrä vaihdella huomattavasti.
- 15 Puheentunnistuksen luotettavuuden parantamiseksi erityisesti vaihtelevissa olosuhteissa on keksinnön erään toisen edullisen suoritusmuodon mukaiseen menetelmään lisätty adaptiivisuutta kynnysarvon θ laskentaan. Tämän adaptiivisuuden aikaansaamiseksi käytetään muutokset $UPDATE_C$, jonka arvo on edullisesti suurempi kuin nolla ja pienempi kuin yksi. Muutokset määritetään aluksi jokin alkuarvo mainitulta arvoalueelta. Tätä muutokset päivitetään puheentunnistuksen aikana edullisesti seuraavasti. Alikeistoista puskureihin tallennettujen näytteiden perusteella lasketaan suurin tehotaso
- 20 win_max ja pienin tehotaso win_min . Tämän jälkeen suoritetaan mainitun lasketun suurimman tehotason win_max vertailu sen hetkiseen tehomaksimiin p_max ja mainitun lasketun pienimmän tehotason win_min vertailu tehominimiin p_min . Jos lasketun suurimman tehotason win_max ja tehomaksimin p_max välisen eron itseisarvo tai tehominimin p_min ja mainitun lasketun pienimmän tehotason win_min välisen eron itseisarvo on kasvanut edellisestä laskentakerrasta, kasvatetaan muutokset $UPDATE_C$. Vastaavasti jos lasketun suurimman tehotason win_max ja tehomaksimin p_max välisen eron itseisarvo tai tehominimin p_min ja mainitun lasketun pienimmän tehotason win_min
- 30 välisen eron itseisarvo on pienentynyt edellisestä laskentakerrasta, pienennetään muutokset $UPDATE_C$. Tämän jälkeen lasketaan uusi tehomaksimi ja tehominimi seuraavasti:
- 35

$$p_min(t) = (1 - UPDATE_C) \cdot p_min(t-1) + (UPDATE_C \cdot win_min)$$

$$p_max(t) = (1 - UPDATE_C) \cdot p_max(t-1) + (UPDATE_C \cdot win_max)$$

- 5 Laskettuja uusia tehomaksimi- ja tehominimiarvoja käytetään seuraavalla näytteenottokierroksella mm. funktion $f()$ suorituksen yhteydessä. Tämän adaptiivisen kertoimen määrittämisen etuna on mm. se, että ympäristöolosuhteissa tapahtuvat muutokset voidaan paremmin huomioida puheentunnistuksessa ja tauon ilmaisu saadaan luotettavammaksi.
- 10 Edellä esitetyt eri toiminnot puheessa olevan tauon ilmaisemiseksi voidaan suurelta osin toteuttaa puheentunnistustilteen kontrollerin ja/tai digitaalisen signaalinkäsittelylaitteen sovellusohjelmistossa. Keksinnön mukaisessa puheentunnistustilteessä voidaan osa toiminnoista, kuten alikaistoihin jako toteuttaa myös analogiatekniikalla, kuten on sinänsä
- 15 tunnettua. Menetelmän suorituksen yhteydessä voidaan eri vaiheissa muodostettavien laskentatulosten, muuttujien jne. tallennuksessa käyttää puheentunnistustilteen muistivälineitä 14, edullisesti luku/kirjoitusmuistia (RAM, Random Access Memory), haihtumatonta, uudelleen kirjoitettavissa olevaa lukumuistia (NVRAM, Non-Volatile RAM),
- 20 FLASH-muistia jne. Myös langattoman viestimen muistivälineitä 22 voidaan käyttää tietojen tallennuksessa.
- Kuvassa 2 keksinnön edullisen suoritusmuodon mukaisesta langattomasta viestimestä MS on esitetty vielä sinänsä tunnetut näppäimistö
- 25 17, näyttölaite 18, digitaalinen/analogiamuunnin 19, kuulokevahvistin 20a, kuuloke 21a, kaiutintoiminnon 2 kuulokevahvistin 20b, kuuloke 21b sekä suurtaajuuslohko 23.
- Nyt esillä olevaa keksintöä voidaan soveltaa useiden eri periaatteella toimivien puheentunnistusjärjestelmien yhteydessä. Keksintö parantaa puheessa olevien taukokohtien ilmaisuvarmuutta, mikä varmentaa varsinaisen puheentunnistuksen tunnistusvarmuutta. Keksinnön mukaista menetelmää käytettäessä ei puheentunnistuksesta ole tarve suorittaa kiinteään aikaikkunaan sidottuna, joten tunnistusviive ei olennaisesti riipu
- 30 siitä, kuinka nopeasti käyttäjä lausuu puhekomentoja. Myös taustamelun vaikutus puheentunnistukseen saadaan keksinnön mukaista menetelmää sovellettaessa pienemmäksi kuin tunnetun tekniikan mukaisissa puheentunnistustilteissä on mahdollista.

On selvää, että keksintöä ei ole rajoitettu ainoastaan edellä esitettyihin suoritusmuotoihin, vaan sitä voidaan muunnella oheisten patenttivaatimusten puitteissa.

Patenttivaatimukset:

1. Menetelmä puheentunnistuksessa puheessa olevien taukojen ilmai-
semiseksi, jossa menetelmässä käyttäjän lausumien puhekomentojen
5 tunnistamiseksi ääni muunnetaan sähköiseksi signaaliksi, **tunnettu**
siitä, että menetelmässä sähköisen signaalin taajuusspektri jaetaan
kahdeksi tai useammaksi alikaistaksi, tallennetaan alikaistojen signaa-
leista näytteitä väliajoin, määritetään alikaistojen energiatasot tallennet-
tujen näytteiden perusteella, määritetään tehon kynnysarvo (thr), ja ver-
10 rataan alikaistojen energiatasoja mainittuun tehon kynnysarvoon (thr),
jolloin vertailutuloksia käytetään tauon ilmaisutuloksen muodostuk-
sessa.
2. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että
15 määritetään ilmaisuaikaraja (END) ja ilmaisumäärä (SB_SUFF_TH),
jolloin menetelmässä alikanavan tauon pituuden laskenta aloitetaan
alikaistan energiatason alittaessa mainitun tehon kynnysarvon (thr),
jolloin menetelmässä muodostetaan alikanavakohtainen ilmaisu las-
kennan saavuttaessa ilmaisuaikarajan (END), tutkitaan, kuinka monel-
20 la alikaistalla energiataso on ollut tehon kynnysarvon (thr) alapuolella
pidempään kuin ilmaisuaikaraja (END), jolloin tauon ilmaisupäätös teh-
dään, jos alikanavakohtaisten ilmaisujen lukumäärä on suurempi tai
yhtä suuri kuin ilmaisumäärä (SB_SUFF_TH).
- 25 3. Patenttivaatimuksen 2 mukainen menetelmä, **tunnettu** siitä, että
menetelmässä lisäksi määritetään aktiivisuusaikaraja
(SB_ACTIVE_TH) ja aktiivisuusmäärä (SB_MIN_TH), jolloin tauon il-
maisupäätös tehdään, jos alikanavakohtaisten ilmaisujen lukumäärä on
suurempi tai yhtäsuuri kuin aktiivisuusmäärä (SB_MIN_TH), ja muilla
30 alikanavilla alikanavan tauon pituuden laskennassa ei ole saavutettu
aktiivisuusaikarajaa (SB_ACTIVE_TH).
4. Patenttivaatimuksen 1, 2 tai 3 mukainen menetelmä, **tunnettu** siitä,
että tehon kynnysarvo (thr) lasketaan kaavalla
35
$$thr = p_{min} + k \cdot (p_{max} - p_{min}),$$
 jossa

p_{\min} = alikanavien tallennetusta näytteistä määritetty pienin
tehomaksimi, ja
 p_{\max} = alikanavien tallennetuista näytteistä määritetty suurin
tehominimi.

5

5. Jonkin patenttivaatimuksen 1—4 mukainen menetelmä, **tunnettu** siitä, että mainittu tehon kynnysarvo (thr) lasketaan adaptiivisesti huomioimalla kulloinenkin ympäristön häiriöäänitaso.

- 10 6. Patenttivaatimuksen 5 mukainen menetelmä, **tunnettu** siitä, että mainitun tehon kynnysarvon (thr) laskemiseksi väliajoin (t) määritetään muutoskerroin (UPDATE_C), ja tallennettujen näytteiden perusteella lasketaan alikaistojen suurin tehotaso (win_max) ja pienin tehotaso (win_min), jolloin määritetään tehomaksimi (p_{\max}) ja tehominimi (p_{\min}) kaavoilla:

15

$$p_{\max}(i,t) = (1 - \text{UPDATE_C}) \cdot p_{\max}(i,t-1) + (\text{UPDATE_C} \cdot \text{win_max})$$

$$p_{\min}(i,t) = (1 - \text{UPDATE_C}) \cdot p_{\min}(i,t-1) + (\text{UPDATE_C} \cdot \text{win_min})$$

- 20 jossa $0 < \text{UPDATE_C} < 1$,
 $0 < i < L$, ja
L on alikaistojen lukumäärä

- 25 7. Patenttivaatimuksen 6 mukainen menetelmä, **tunnettu** siitä, että menetelmässä lisäksi:

- kasvatetaan muutoskerrointa (UPDATE_C), mikäli mainitun lasketun suurimman tehotason (win_max) ja tehomaksimin (p_{\max}) välisen eron itseisarvo tai tehominimin (p_{\min}) ja mainitun lasketun pienimmän tehotason (win_min) välisen eron itseisarvo on kasvanut,
- 30 – pienennetään muutoskerrointa (UPDATE_C), mikäli mainitun lasketun suurimman tehotason (win_max) ja tehomaksimin (p_{\max}) välisen eron itseisarvo tai tehominimin (p_{\min}) ja mainitun lasketun pienimmän tehotason (win_min) välisen eron itseisarvo on pienentynyt.

35

8. Puheentunnistuslaite (16), joka käsittää välineet (1a, 1b) käyttäjän lausumien puhekomentojen muuntamiseksi sähköiseksi signaaliksi, **tunnettu** siitä, että se käsittää lisäksi:

- 5 – välineet (8) sähköisen signaalin taajuusspektrin jakamiseksi kahdeksi tai useammaksi alikaistaksi,
- välineet (14) näytteiden tallentamiseksi väliajoin alikaistojen signaaleista,
- välineet (5, 13) energiatasojen määrittämiseksi alikaistoista tallennettujen näytteiden perusteella,
- 10 – välineet (5, 13) tehon kynnysarvon (thr) määrittämiseksi,
- välineet (5, 13) alikaistojen energiatasojen vertailemiseksi mainittuun tehon kynnysarvoon (thr), ja
- välineet (5, 13) puheessa olevan tauon ilmaisemiseksi mainittujen vertailutulosten perusteella.

15

9. Patenttivaatimuksen 8 mukainen puheentunnistuslaite (16), **tunnettu** siitä, että tehon kynnysarvo (thr) on laskettu kaavalla

$$thr = p_{min} + k \cdot (p_{max} - p_{min}), \text{ jossa}$$

20

p_{min} = alikanavien tallennetuista näytteistä määritetty pienin tehomaksimi, ja
 p_{max} = alikanavien tallennetuista näytteistä määritetty suurin tehominimi.

25

10. Patenttivaatimuksen 8 tai 9 mukainen puheentunnistuslaite (16), **tunnettu** siitä, että se käsittää lisäksi välineet (10, 11) alikaistojen signaalien suodattamiseksi ennen tallennusta.

30

11. Langaton viestin (MS), joka käsittää välineet (16) puheen tunnistamiseksi, ja välineet (1a, 1b) käyttäjän lausumien puhekomentojen muuntamiseksi sähköiseksi signaaliksi, **tunnettu** siitä, että välineet (16) puheen tunnistamiseksi käsittää lisäksi:

35

- välineet (8) sähköisen signaalin taajuusspektrin jakamiseksi kahdeksi tai useammaksi alikaistaksi,
- välineet (14) näytteiden tallentamiseksi väliajoin alikaistojen signaaleista,

- välineet (5, 13) energiatasojen määrittämiseksi alikaistoista tallennettujen näytteiden perusteella,
 - välineet (5, 13) tehon kynnysarvon (thr) määrittämiseksi,
 - välineet (5, 13) alikaistojen energiatasojen vertailemiseksi mainittuun tehon kynnysarvoon (thr), ja
 - välineet (5, 13) puheessa olevan tauon ilmaisemiseksi mainittujen vertailutulosten perusteella.
- 5

18

L 3

(57) Tiivistelmä

Menetelmässä puheessa olevien taukojen ilmaisemiseksi käyttäjän lausumien puhekomentojen tunnistamista varten ääni muunnetaan sähköiseksi signaaliksi, jonka taajuusspektri jaetaan kahdeksi tai useammaksi alikaistaksi. Aikaistojen signaaleista tallennetaan näytteitä väliajoin, määritetään aikaistojen energiatasot tallennettujen näytteiden perusteella, määritetään tehon kynnysarvo (thr), ja verrataan aikaistojen energiatasoa mainittuun tehon kynnysarvoon (thr). Vertailutuloksia käytetään tauon ilmaisutuloksen muodostuksessa.

5 Fig. 1

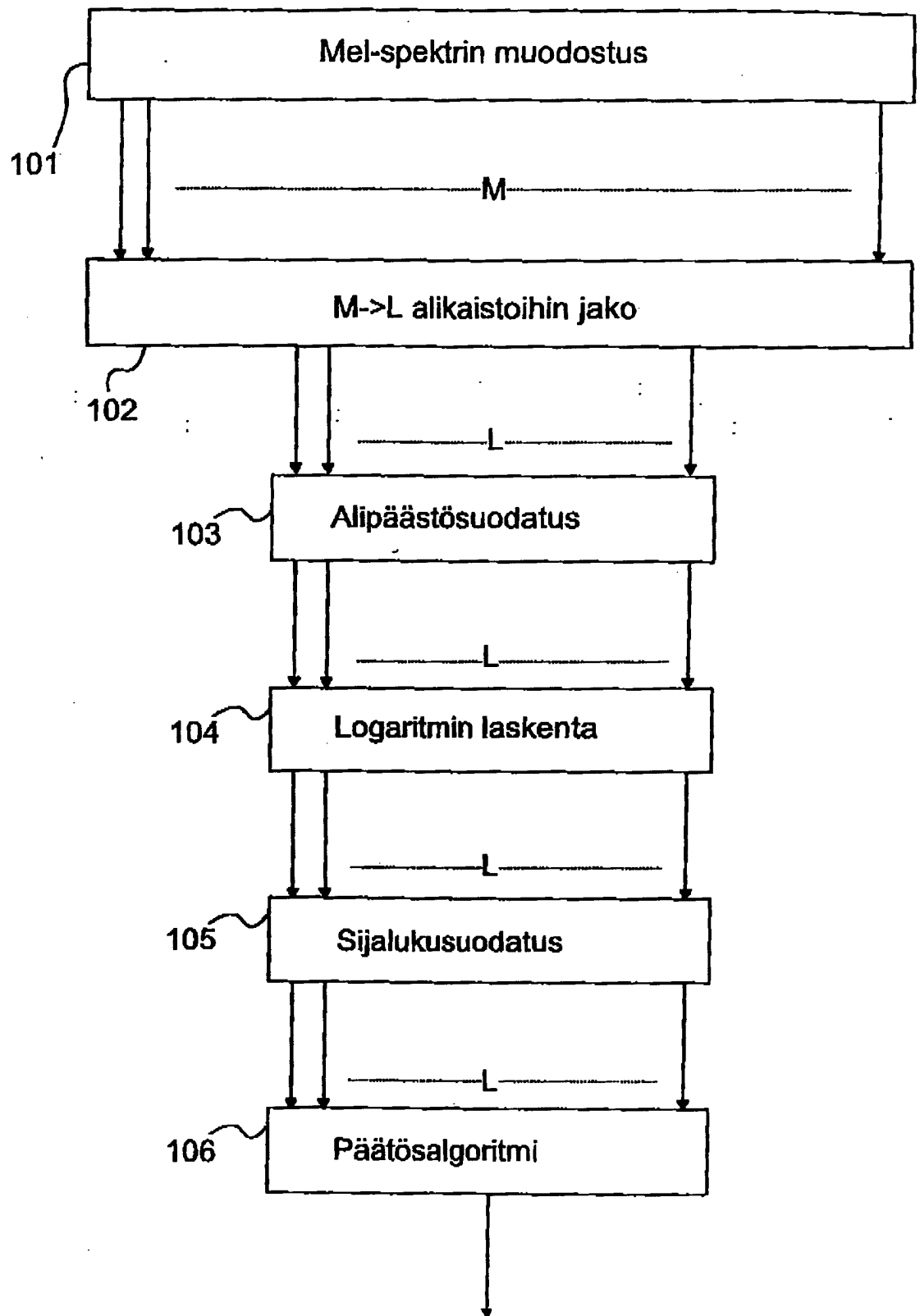


Fig 1

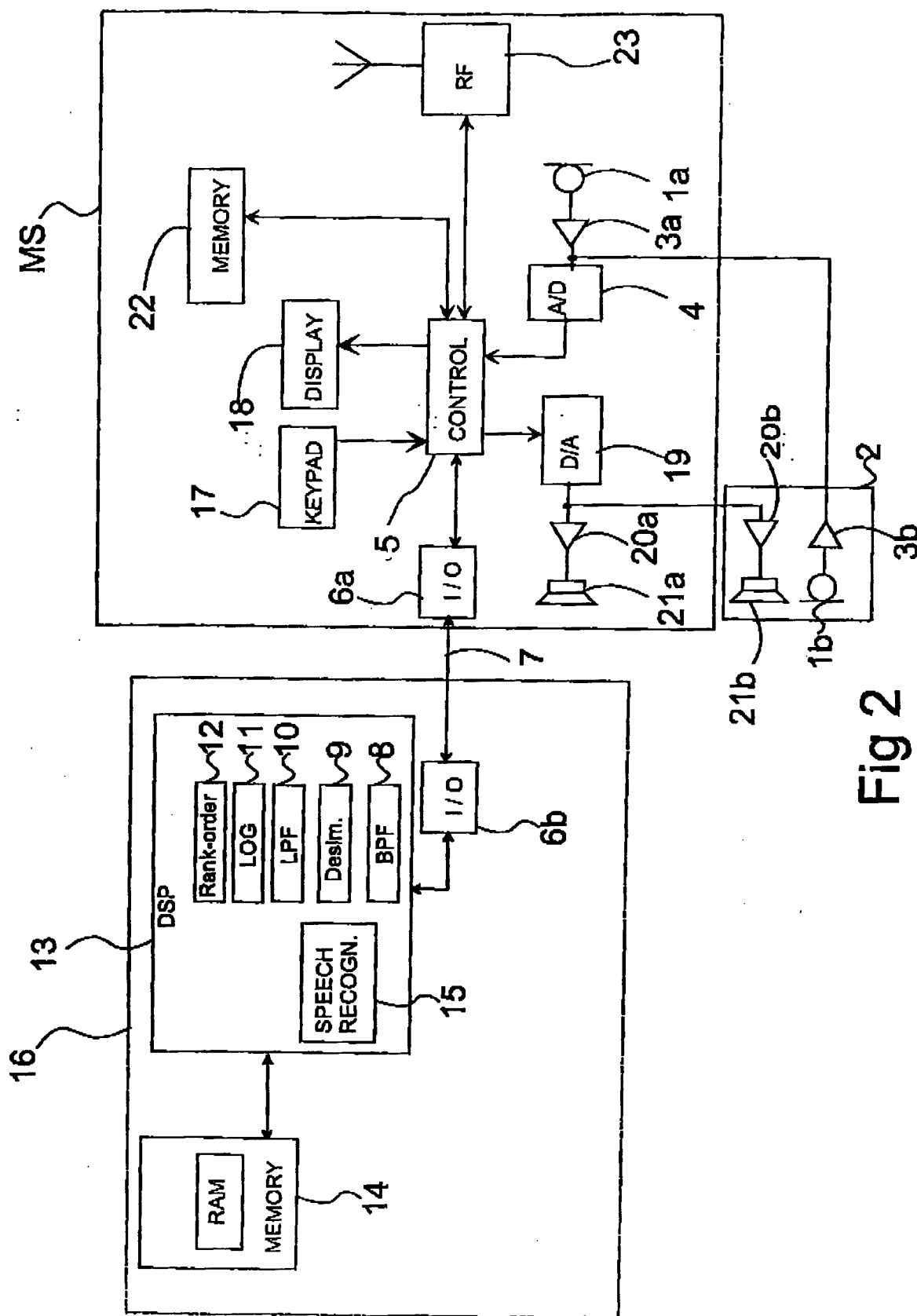


Fig 2

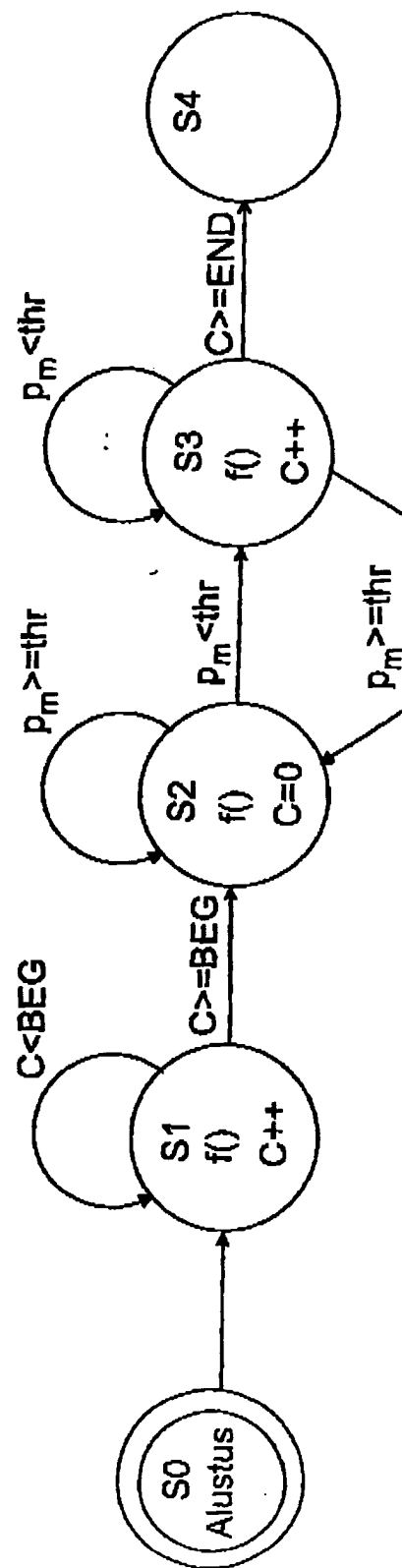


Fig 3

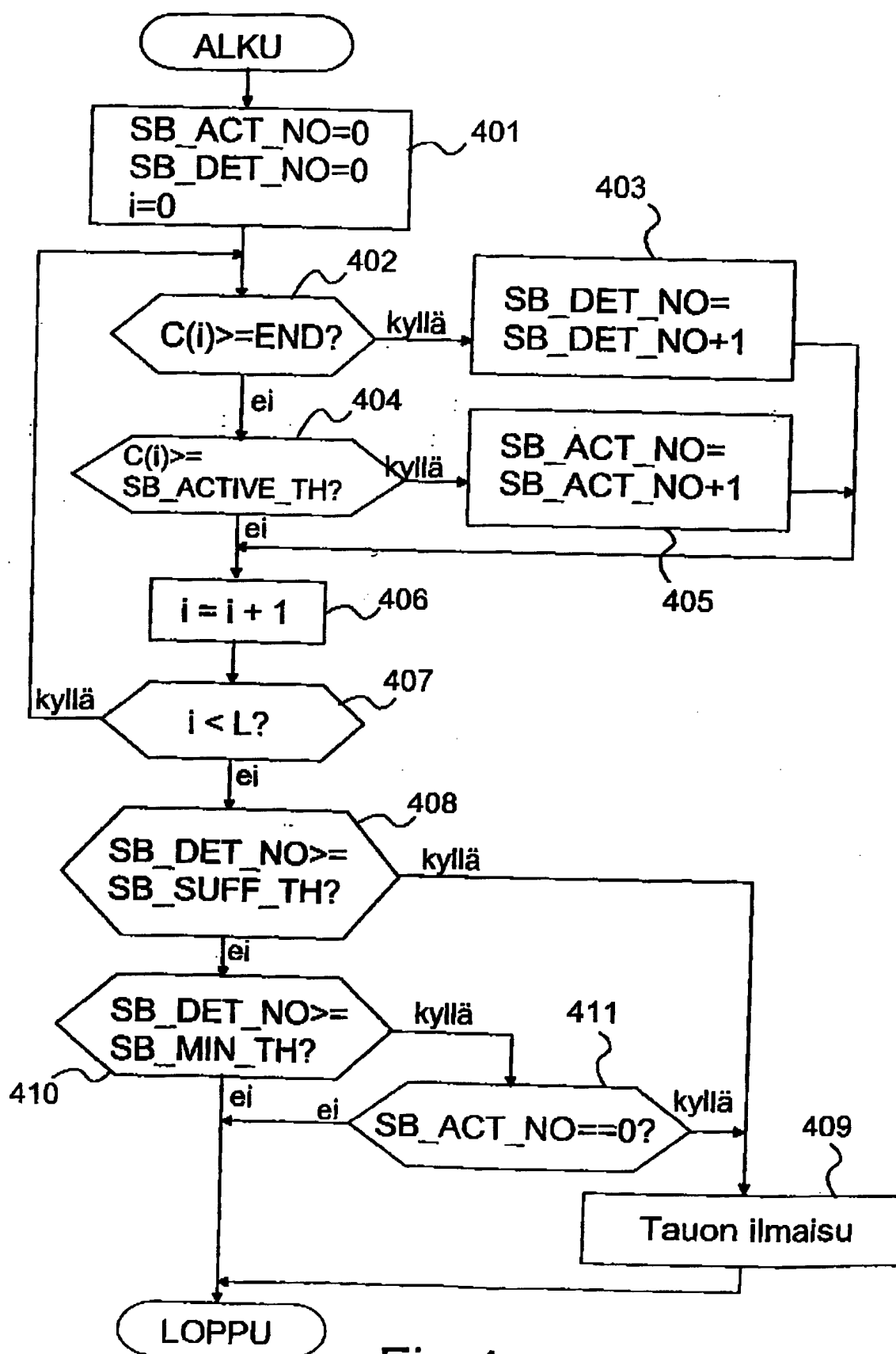


Fig 4